

Multivariable Analysis

**A Practical Guide for
Clinicians**

MITCHELL H. KATZ



PUBLISHED BY THE PRESS SYNDICATE OF THE UNIVERSITY OF CAMBRIDGE
The Pitt Building, Trumpington Street, Cambridge, United Kingdom

CAMBRIDGE UNIVERSITY PRESS

The Edinburgh Building, Cambridge CB2 2RU, UK <http://www.cup.cam.ac.uk>
40 West 20th Street, New York, NY 10011-4211, USA <http://www.cup.org>
10 Stamford Road, Oakleigh, Melbourne 3166, Australia

© Cambridge University Press 1999

This book is in copyright. Subject to statutory exception
and to the provisions of relevant collective licensing agreements,
no reproduction of any part may take place without
the written permission of Cambridge University Press.

First published 1999

Printed in the United States of America

Typeset in Stone Serif 9.5/13pt. and Avenir in $\text{\LaTeX} 2_{\epsilon}$ [TB]

*A catalog record for this book is available from
the British Library*

Library of Congress Cataloging-in-Publication Data

Katz, Mitchell H.

Multivariable analysis : a practical guide for clinicians /

Mitchell H. Katz.

p. cm.

1. Medicine—Research—Statistical methods. 2. Multivariate
analysis. 3. Biometry. 4. Medical statistics. I. Title.

R853.S7K38 1999

610'.7'27—dc21

98-39350

CIP

ISBN 0 521 59301 8 hardback

ISBN 0 521 59693 9 paperback

Contents

<i>Preface</i>	<i>page xiii</i>
1 Introduction	1
1.1 Why should I do multivariable analysis? • 1	
1.2 What are confounders and how does multivariable analysis help me to deal with them? • 6	
1.3 What are suppressers and how does multivariable analysis help me to deal with them? • 11	
1.4 What are interactions and how does multivariable analysis help me to deal with them? • 13	
2 Common Uses of Multivariable Models	17
2.1 What are the most common uses of multivariable models in clinical research? • 17	
2.2 How do I choose what type of multivariable analysis to use? • 26	
3 Outcome Variables in Multivariable Analysis	27
3.1 How does the nature of my outcome variable influence my choice of which type of multivariable analysis to do? • 27	
3.2 What should I do if my outcome variable is ordinal or nominal? • 27	
3.3 What are the advantages of using time to occurrence of a dichotomous event instead of the simpler cumulative outcome of a dichotomous event at a point in time? • 29	

4	Independent Variables in Multivariable Analysis	33
4.1	What kind of independent variables can I use with multivariable analyses? • 33	
4.2	What should I do with my ordinal and nominal independent variables? • 33	
5	Assumptions of Multiple Linear Regression, Logistic Regression, and Proportional Hazards Analysis	36
5.1	What are the assumptions of multiple linear regression, multiple logistic regression, and proportional hazards analysis? • 36	
5.2	What is being modeled in multiple linear regression, multiple logistic regression, and proportional hazards analysis? • 36	
5.3	What is the relationship of multiple independent variables to outcome in multiple linear regression, multiple logistic regression, and proportional hazards analysis? • 41	
5.4	What is the relationship of an interval-independent variable to the outcome in multiple linear regression, multiple logistic regression, and proportional hazards analysis? • 41	
5.5	What if my interval-independent variable does not have a linear relationship with my outcome? • 45	
5.6	Assuming that my interval-independent variable fits a linear assumption, is there any reason to group it into interval categories or create multiple dichotomous variables? • 50	
5.7	What are the assumptions about the distribution of the outcome and the variance? • 51	
5.8	What should I do if I find significant violations of the assumptions of normal distribution and equal variance in my multiple linear regression analysis? • 54	
6	Relationship of Independent Variables to One Another	55
6.1	Does it matter if my independent variables are related to each other? • 55	

6.2	How do I assess whether my variables are multicollinear? • 56	
6.3	What should I do with multicollinear variables? • 58	
7	Setting Up a Multivariable Analysis: Subjects	60
7.1	How many subjects do I need to do multivariable analyses? • 60	
7.2	What if I have too many independent variables given my sample size? • 64	
7.3	What if some of my subjects do not complete my study? • 70	
7.4	What assumptions does censoring make about observations with unequal lengths of follow-up? • 71	
7.5	How likely is it that the censoring assumption is valid in my study? • 75	
7.6	How can I test the validity of the censoring assumption for my data? • 80	
8	Performing the Analysis	84
8.1	What numbers should I assign for dichotomous or ordinal variables in my analysis? • 84	
8.2	Does it matter what I choose as my reference category for multiple dichotomous (“dummied”) variables? • 85	
8.3	How do I enter interaction terms into my analysis? • 87	
8.4	How do I enter time into my proportional hazards or other survival analysis? • 89	
8.5	What about subjects who experience their outcome on their start date? • 95	
8.6	What about subjects who have a survival time shorter than physiologically possible? • 97	
8.7	What should I do about missing data on my independent variables? • 99	
8.8	What should I do about missing data on my outcome variable? • 108	
8.9	What are variable selection techniques? Which variable selection technique should I use? • 110	

- 8.10** If I use a forward or backward selection technique, what level of statistical significance should I set for inclusion/exclusion of a variable? • 114
- 8.11** Do I have to use a variable selection technique at all? • 115
- 8.12** What value should I specify for tolerance in my logistic regression or proportional hazards model? • 115
- 8.13** How many iterations (attempts to solve) should I specify for my logistic regression or proportional hazards model? • 116
- 8.14** What value should I specify for the convergence criteria for my logistic regression or proportional hazards model? • 116
- 8.15** My model will not converge. What should I do? • 116

9 Interpreting the Analysis

118

- 9.1** What information will the printout from my analysis provide? • 118
- 9.2** How do I assess how well my model accounts for my outcome? • 118
- 9.3** What do the coefficients tell me about the relationship between each variable and the outcome? • 127
- 9.4** How do I get odds ratios and relative hazards from the multivariable analysis? What do they mean? • 128
- 9.5** How do I interpret the odds ratio and relative hazard when the independent variable is interval? • 132
- 9.6** How do I compute the confidence intervals for the odds ratios and relative hazards? • 133
- 9.7** What are standardized coefficients and should I use them? • 134
- 9.8** How do I test the statistical significance of my coefficients? • 135
- 9.9** How do I interpret the results of interaction terms? • 138

9.10	Do I have to adjust my multivariable regression coefficients for multiple comparisons? • 138	
10	Checking the Assumptions of the Analysis	141
10.1	How do I know if my data fit the assumptions of my multivariable model? • 141	
10.2	How do I assess the linearity, normal distribution, and equal variance assumptions of multiple linear regression? • 142	
10.3	How do I assess the linearity assumption of multiple logistic regression and proportional hazards analysis? • 143	
10.4	What are outliers and how do I detect them in my multiple linear regression models? • 144	
10.5	How do I detect outliers in my multiple logistic regression model? • 146	
10.6	What about analysis of residuals with proportional hazards analysis? • 146	
10.7	What should I do when I detect outliers? • 146	
10.8	What is the additive assumption and how do I assess whether my multiple independent variables fit this assumption? • 147	
10.9	What does the addition assumption mean for interval-independent variables • 150	
10.10	What is the proportionality assumption? • 151	
10.11	How do I test the proportionality assumption? • 153	
10.12	What if the proportionality assumption does not hold for my data? • 156	
11	Validation of Models	158
11.1	How can I validate my models? • 158	
12	Special Topics	164
12.1	What if my data set has matched cases and controls? • 164	
12.2	What if my data set has repeated observations of outcome for the same individuals? • 166	
12.3	What if my outcome can occur in more than one body part in the same person? • 170	

12.4	What if the independent variable changes value during the course of the study? • 172	
12.5	What are the advantages and disadvantages of time-dependent covariates? • 173	
12.6	What if the frequency of my outcome is really low over time (rare disease)? • 175	
12.7	What are classification and regression trees (CART) and should I use them? • 176	
12.8	How can I get best use of my biostatistician? • 179	
12.9	How do I choose which software package to use? • 180	
13	Publishing Your Study	181
13.1	How much information about how I constructed my multivariable models should I put in the Methods section? • 181	
13.2	Do I need to cite a statistical reference for my choice of multivariable models? • 183	
13.3	Which parts of my multivariable analysis should I report in the Results section? • 183	
14	Summary: Steps for Constructing a Multivariable Model	187

1

Introduction

1.1 Why should I do multivariable analysis?

We live in a multivariable world. Most events, whether medical, political, social, or personal, have multiple causes. And these causes are related to one another. Multivariable analysis¹ is a statistical tool for determining the relative contributions of different causes to a single event or outcome.

Clinical researchers, in particular, need multivariable analysis because most diseases have multiple causes and prognosis is usually determined by a large number of factors. Even for those infectious diseases that are known to be caused by a single pathogen, a number of factors affect whether an exposed individual becomes ill, including the characteristics of the pathogen (e.g., virulence of strain), the route of exposure (e.g., respiratory route), the intensity of exposure (e.g., size of inoculum), and the host response (e.g., immunologic defense).

Multivariable analysis allows us to sort out the multifaceted nature of risk factors and their relative contribution to outcome. For example, observational epidemiology has taught us that there are a number of risk factors associated with premature mortality, notably smoking, a sedentary lifestyle, obesity, elevated cholesterol, and hypertension. Note that I did not say that these factors *cause* premature mortality. Statistics alone cannot prove that a relationship between a risk factor

DEFINITION

Multivariable analysis is a tool for determining the relative contributions of different causes to a single event.

¹ The terms “multivariate analysis” and “multivariable analysis” are often used interchangeably. In the strict sense, multivariate analysis refers to simultaneously predicting multiple outcomes. Since this book deals with techniques that use multiple variables to predict a single outcome, I prefer the more general term multivariable analysis.

and an outcome are causal.² Causality is established on the basis of biological plausibility and rigorous study designs, such as randomized controlled trials, which eliminate sources of potential bias.

Identification of risk factors of premature mortality through observational studies has been particularly important because you cannot randomize people to many of the conditions that cause premature mortality, such as smoking, sedentary lifestyle, or obesity. And yet these conditions tend to occur together; that is, people who smoke tend to exercise less and be more likely to be obese. How does multi-variable analysis separate the *independent* contribution of each of these factors? Let's consider the case of exercise. Numerous studies have shown that persons who exercise live longer than persons with sedentary lifestyles. But if the only reason that persons who exercise live longer is that they are less likely to smoke and more likely to eat low fat meals leading to lower cholesterol, then initiating an exercising routine would not change a person's life expectancy.

The Aerobics Center Longitudinal Study tackled this important question.³ They evaluated the relationship between exercise and mortality in 25,341 men and 7,080 women. All participants had a baseline examination between 1970 and 1989. The examination included a physical examination, laboratory tests, and a treadmill evaluation to assess physical fitness. Participants were followed for an average of 8.4 years for the men and 7.5 years for the women.

Table 1.1 compares the characteristics of survivors to persons who had died during the follow-up. You can see that there are a number of significant differences between survivors and decedents among men and women. Specifically, survivors were younger, had lower blood pressure, lower cholesterol, were less likely to smoke, and were more physically fit (based on the length of time they stayed on the treadmill and their level of effort).

Although the results are interesting, Table 1.1 does not answer our basic question: Does being physically fit independently increase

² Throughout the text I use the terms "associated with" and "related to" interchangeably. Similarly, I use the terms "risk factor" and "independent variable," and the terms "outcome" and "dependent variable," interchangeably. Although many use the term "predicts" to refer to the association between an independent variable and an outcome, the term implies causality and I prefer to reserve it for when we are determining how well a model predicts the outcome of individual subjects (Section 9.2C).

³ Blair, S.N., Kampert, J.B., Kohl, H.W., et al. "Influences of cardiorespiratory fitness and other precursors on cardiovascular disease and all-cause mortality in men and women." *JAMA* 1996;276:205–10.

TABLE 1.1

Baseline characteristics of survivors and decedents, Aerobics Center Longitudinal Study.

Characteristics	Men		Women	
	Survivors (n = 24,740)	Decedents (n = 601)	Survivors (n = 6,991)	Decedents (n = 89)
Age, y (SD)	42.7 (9.7)	52.1 (11.4)	42.6 (10.9)	53.3 (11.2)
Body mass index, kg/m ² (SD)	26.0 (3.6)	26.3 (3.5)	22.6 (3.9)	23.7 (4.5)
Systolic blood pressure, mm Hg (SD)	121.1 (13.5)	130.4 (19.1)	112.6 (14.8)	122.6 (17.3)
Total cholesterol, mg/dL (SD)	213.1 (40.6)	228.9 (45.4)	202.7 (40.5)	228.2 (40.8)
Fasting glucose, mg/dL (SD)	100.4 (16.3)	108.1 (32.0)	94.4 (14.5)	99.9 (25.0)
Fitness, %				
Low	20.1	41.6	18.8	44.9
Moderate	42.0	39.1	40.6	33.7
High	37.9	19.3	40.6	21.3
Current or recent smoker, %	26.3	36.9	18.5	30.3
Family history of coronary heart disease, %	25.4	33.8	25.2	27.0
Abnormal electrocardiogram, %	6.9	26.3	4.8	18.0
Chronic illness, %	18.4	40.3	13.4	20.2

Adapted with permission from Blair, S.N., et al. "Influences of cardiorespiratory fitness and other precursors on cardiovascular disease and all-cause mortality in men and women." *JAMA* 1996;276:205–10. Copyright 1996, American Medical Association. Additional data provided by authors.

longevity? It doesn't answer the question because whereas the high-fitness group was less likely to die during the study period, those who were physically fit may just have been younger, been less likely to smoke, or had lower blood pressure.

To determine whether exercise is independently associated with mortality, the authors performed proportional hazards analysis, a type of multivariable analysis. The results are shown in Table 1.2. If you compare the number of deaths per thousand person-years in men, you can see that there were more deaths in the low-fitness group (38.1)

TABLE 1.2

Multivariable analysis of risk factors for all-cause mortality, Aerobics Center Longitudinal Study.

Independent variable	Men		Women	
	Deaths per 10,000 person-years	Adjusted relative risk (95% CI)	Deaths per 10,000 person-years	Adjusted relative risk (95% CI)
Fitness				
Low	38.1	1.52 (1.28–1.82)	27.8	2.10 (1.36–3.26)
Moderate/High	25.0	1.0 (ref.)	13.2	1.0 (ref.)
Smoking Status				
Current or recent smoker	39.4	1.65 (1.39–1.97)	27.8	1.99 (1.25–3.17)
Past or never smoked	23.9	1.0 (ref.)	14.0	1.0 (ref.)
Systolic blood pressure				
≥140 mm Hg	35.6	1.30 (1.08–1.58)	13.0	0.76 (0.41–1.40)
<140 mm Hg	27.3	1.0 (ref.)	17.1	1.0 (ref.)
Cholesterol				
≥240 mg/dL	35.1	1.34 (1.13–1.59)	18.0	1.09 (0.68–1.74)
<240 mg/dL	26.1	1.0 (ref.)	16.6	1.0 (ref.)
Family history of coronary heart disease				
Yes	29.9	1.07 (0.90–1.29)	12.8	0.70 (0.43–1.16)
No	27.8	1.0 (ref.)	18.2	1.0 (ref.)
Body mass index				
≥27 kg/m ²	28.8	1.02 (0.86–1.22)	15.9	0.94 (0.52–1.69)
<27 kg/m ²	28.2	1.0 (ref.)	16.9	1.0 (ref.)
Fasting glucose				
≥120 mg/dL	34.4	1.24 (0.98–1.56)	29.6	1.79 (0.80–4.00)
<120 mg/dL	27.9	1.0 (ref.)	16.5	1.0 (ref.)
Abnormal electrocardiogram				
Yes	44.4	1.64 (1.34–2.01)	25.3	1.55 (0.87–2.77)
No	27.1	1.0 (ref.)	16.3	1.0 (ref.)
Chronic illness				
Yes	41.2	1.63 (1.37–1.95)	17.5	1.05 (0.61–1.82)
No	25.3	1.0 (ref.)	16.7	1.0 (ref.)

Adapted with permission from Blair, S.N., et al. "Influences of cardiorespiratory fitness and other precursors on cardiovascular disease and all-cause mortality in men and women." *JAMA* 1996;276:205–10. Copyright 1996, American Medical Association. Additional data provided by authors.

TABLE 1.3

Stratified analysis of smoking and fitness on all-cause mortality among men, Aerobics Center Longitudinal Study.

	Deaths per 10,000 person-years	Stratum-specific relative risk (95% CI)
Smokers		
Low fitness	48.0	1.63 (1.26–2.13)
Moderate/high fitness	29.4	1.0 (ref.)
Nonsmokers		
Low fitness	44.0	2.19 (1.77–2.70)
Moderate/high fitness	20.1	1.0 (ref.)

Data supplied by Aerobics Center Longitudinal Study.

than in the moderate/high-fitness group (25.0). This difference is reflected in the elevated relative risk for lower fitness ($38.1/25.0 = 1.52$). These results are adjusted for all of the other variables listed in the table. This means that low fitness is associated with higher mortality independent of the effects of other known risk factors for mortality, such as smoking, elevated blood pressure, cholesterol, and family history. A similar pattern is seen for women.

Was there any way to answer this question without multivariable analysis? One could have performed stratified analysis. Stratified analysis assesses the effect of a risk factor on outcome while holding another variable constant. So, for example, we could compare physically fit to unfit persons separately among smokers and nonsmokers. This would allow us to calculate a relative risk for the impact of fitness on mortality, independent of smoking. This analysis is shown in Table 1.3.

Unlike the multivariable analysis in Table 1.2, the analyses in Table 1.3 are bivariate.⁴ We see that the mortality rate is greater among those at low fitness compared to those at moderate/high fitness both among smokers (48.0 vs. 29.4) and among nonsmokers (44.0 vs. 20.1).

⁴ Some researchers use the term “univariate” to describe the association between two variables. I think it is more informative to restrict the term univariate to analyses of a single variable (e.g., mean, median), while using the term “bivariate” to refer to the association between two variables.

DEFINITION

Stratified analysis assesses the effect of a risk factor on outcome while holding another variable constant.

This stratified analysis shows that the effect of fitness is independent of smoking status.

But what about all of the other variables that might affect the relationship between fitness and longevity? You could certainly stratify for each one individually, proving that the effect of fitness on longevity is independent not only of smoking status, but also independent of elevated cholesterol, elevated blood pressure, and so on. However, this would only prove that the relationship is independent of these variables taken singly. To stratify by two variables (smoking and cholesterol), you would have to assess the relationship between fitness and mortality in four groups (smokers with high cholesterol; smokers with low cholesterol; nonsmokers with high cholesterol; nonsmokers with low cholesterol). To stratify by three variables (smoking status, cholesterol level, and elevated blood pressure (yes/no)), you would have to assess the relationship between fitness and mortality in eight groups; add elevated glucose (yes/no) and you would have 16 groups; add age (in six decades) and you would have 96 groups; and we haven't even yet taken into account all of the variables in Table 1.1 that are associated with mortality.

With each stratification variable you add, you increase the number of subgroups for which you have to individually assess whether the relationship between fitness and mortality holds. Besides producing mountains of printouts, and requiring a book (rather than a journal article) to report your results, you would likely have an insufficient sample size in some of these subgroups, even if you started with a large sample size. For example, in the Aerobics Center Longitudinal Study there were 25,341 men but only 601 deaths. With 96 subgroups, assuming uniform distributions, you would expect only about 6 deaths per subgroup. But, in reality you wouldn't have uniform distributions. Some samples would be very small, and some would have no outcomes at all.

Multivariable analysis overcomes this limitation. It allows you to simultaneously assess the impact of multiple independent variables on outcome. But there is (always) a cost: The model makes certain assumptions about the nature of the data. These assumptions are sometimes hard to verify. We will take up these issues in Chapters 5, 6, 7, and 10.

1.2 What are confounders and how does multivariable analysis help me to deal with them?

The ability of multivariable analysis to *simultaneously* assess the independent contribution of a number of risk factors to outcome is

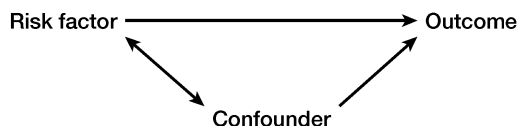


Figure 1.1. Relationships among risk factor, confounder, and outcome.

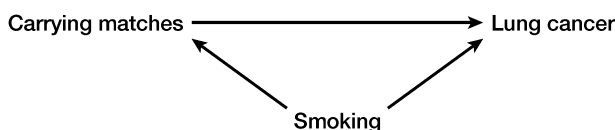


Figure 1.2. Relationships among carrying matches, smoking, and lung cancer.

particularly important when you have confounding. Confounding occurs when the apparent association between a risk factor and an outcome is affected by the relationship of a third variable to the risk factor and the outcome; the third variable is called a confounder.

For a variable to be a confounder, the variable must be associated with the risk factor and causally related to the outcome (Figure 1.1).

A classically taught example of confounding is the relationship between carrying matches and developing lung cancer (Figure 1.2). Persons who carry matches have a greater chance of developing lung cancer; the confounder is smoking. This example is often used to illustrate confounding because it is easy to grasp that carrying matches cannot possibly cause lung cancer.

Stratified analysis can be used to assess and eliminate confounding. If you stratify by smoking status you will find that carrying matches is not associated with lung cancer. That is, there will be no relationship between carrying matches and lung cancer when you look separately among smokers and nonsmokers. The statistical evidence of confounding is the difference between the unstratified and the stratified analysis. In the unstratified analysis the chi-square test would be significant and the odds ratio for the impact of matches on lung cancer would be significantly greater than one. In the two stratified analyses (smokers and nonsmokers), carrying matches would not be significantly associated with lung cancer; the odds ratio would be one in both strata. This differs from the example of stratified analysis in Table 1.3 where exercise was significantly associated with mortality for both smokers and nonsmokers.

Most clinical examples of confounding are more subtle and harder to diagnose than the case of matches and lung cancer. Let's look at the

DEFINITION

A *confounder* is associated with the risk factor and causally related to the outcome.

TABLE 1.4

Bivariate association between smoking status and risk of death.

Bivariate	Nonsmokers	Former smokers	Recent quitters	Persistent smokers
Relative risk of death	1.0 (ref.)	1.08 (.92–1.26)	.56 (.40–.77)	.74 (.59–.94)

Adapted from Hasdai, D., et al. "Effect of smoking status on the long-term outcome after successful percutaneous coronary revascularization." *N. Engl. J. Med.* 1997;336:755–61.

relationship between smoking and prognosis in patients with coronary artery disease following angioplasty (the opening of clogged coronary vessels with the use of a wire and a balloon).

Everyone knows (although the cigarette companies long claimed ignorance) that smoking increases the risk of death. Countless studies including the Aerobics Center Longitudinal Study (Table 1.2) have demonstrated that smoking is associated with increased mortality. How then can we explain the results of Hasdai and colleagues?⁵ They followed 5,437 patients with coronary artery disease who had angioplasty. They divided their sample into nonsmokers, former smokers (quit at least six months before procedure), quitters (quit immediately following the procedure), and persistent smokers. The relative risk of death with the 95% confidence intervals are shown in Table 1.4.

How can the risk of death be lower among persons who persistently smoke than those who never smoked? In the case of recent quitters, you would expect their risk of death to return toward normal only after years of not smoking – and even then you wouldn't actually expect quitters to have a lower risk of death.

Before you assume that there is something wrong with this study, several other studies have found a similar relationship between smoking and better prognosis among patients with coronary artery disease after thrombolytic therapy. This effect has been named the "smoker's

⁵ Hasdai, D., Garratt, K.N., Grill, D.E., Lerman, A., Homes, D.R. "Effect of smoking status on the long-term outcome after successful percutaneous coronary revascularization." *N. Engl. J. Med.* 1997;336:755–61.

TABLE 1.5

Association between demographic and clinical factors and smoking status.

	Nonsmokers	Former smokers	Recent quitters	Persistent smokers
Age, year \pm SD	67 \pm 11	65 \pm 10	56 \pm 10	55 \pm 11
Duration of angina, month \pm SD	41 \pm 66	51 \pm 72	21 \pm 46	29 \pm 55
Diabetes, %	21%	18%	8%	10%
Hypertension, %	54%	48%	38%	39%
Extent of coronary artery disease, %				
One vessel	50%	51%	57%	55%
Two vessels	36%	36%	34%	36%
Three vessels	14%	13%	10%	9%

Adapted from Hasdai, D., et al. "Effect of smoking status on the long-term outcome after successful percutaneous coronary revascularization." *N. Engl. J. Med.* 1997;336:755–61.

paradox."⁶ What is behind the paradox? Look at Table 1.5. As you can see, compared to nonsmokers and former smokers, quitters and persistent smokers are younger, have had angina for a shorter period of time, are less likely to have diabetes and hypertension, and have less severe coronary artery disease (i.e., more one-vessel disease and less three-vessel disease). Given this, it is not so surprising that the recent quitters and persistent smokers have a lower risk of death than nonsmokers and former smokers: They are younger and have fewer underlying medical problems than the nonsmokers and former smokers.

Compare the bivariate (unadjusted) risk of death to the multivariable risk of death (Table 1.6). Note that in the multivariable analysis the researchers adjusted for those differences, such as age and duration of angina, that existed among the four groups.

With statistical adjustment for the baseline differences between the groups, the quitters and persistent smokers have a significantly

✓ TIP

Multivariable analysis is preferable to stratified analysis when you have multiple confounders.

⁶ Barbash, G.I., Reiner, J., White, H.D., et al. "Evaluation of paradoxical beneficial effects of smoking in patients receiving thrombolytic therapy for acute myocardial infarction: Mechanisms of the 'smoker's paradox' from the GUSTO-I trial, with angiographic insights." *J. Am. Coll. Cardiol.* 1995;26:1222–9.

TABLE 1.6

Comparison of bivariate and multivariable association between smoking status and risk of death.

	Nonsmokers	Former smokers	Recent quitters	Persistent smokers
Relative risk of death				
Bivariate	1.0 (ref.)	1.08 (.92–1.26)	.56 (.40–.77)	.74 (.59–.94)
Relative risk of death				
Multivariable	1.0 (ref.)	1.34 (1.14–1.57)	1.21 (.87–1.70)	1.76 (1.37–2.26)

Adapted from Hasdai, D., et al. "Effect of smoking status on the long-term outcome after successful percutaneous coronary revascularization." *N. Engl. J. Med.* 1997;336:755–61.

greater risk of death than nonsmokers – a much more sensible result. (The quitters also have a greater risk of death than the nonsmokers, but the confidence intervals of the relative risk do not exclude one.) The difference between the bivariate and multivariable analysis indicates that confounding is present. The advantage of multivariable analysis over stratified analysis is that it would have been difficult to stratify for age, duration of angina, diabetes, hypertension, and extent of coronary artery disease.

DEFINITION

An *intervening variable* is on the causal pathway to your outcome.

Although the use of multivariable models to adjust for multiple confounders has been a major boon for epidemiology, it is possible to be over zealous in adjusting for potential confounders and thereby adjust away the very effect you are trying to demonstrate. Camargo and colleagues recognized this in their study of the relationship between moderate alcohol consumption and risk of heart attack.⁷ Sensibly, they adjusted for age, smoking, exercise, diabetes, and family history of heart attack. However, they did not adjust for blood pressure, body mass index, or hypercholesterolemia. Why not? After all, these factors fit the definition of a confounder, in that they are associated with the risk factor (alcohol consumption) and causally related to the outcome (myocardial infarction). The problem is that alcohol consumption can cause elevations in blood pressure, body mass index, and hypercholesterolemia. Therefore, as illustrated in Figure 1.3, these variables may be

⁷ Camargo, C.A., Stampfer, M.J., Glynn, R.J., et al. "Moderate alcohol consumption and risk for angina pectoris or myocardial infarction in U.S. male physicians." *Ann. Intern. Med.* 1997;126:372–5.

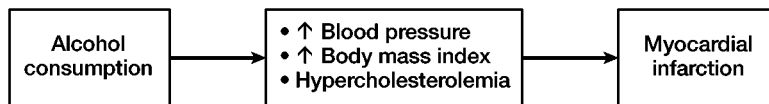


Figure 1.3. Hypothesized pathway by which alcohol consumption may cause myocardial infarction.



Figure 1.4. Relationships among risk factor, suppresser, and outcome.

on the causal pathway to myocardial infarction and should be thought of as intervening variables rather than as confounders. If you adjust for intervening variables, as if they were confounders, you will adjust away the effect you are trying to demonstrate.

Statistics cannot tell you whether something is a confounder or an intervening variable. Statistically, confounders and intervening variables operate the same. Whether to include a variable in your model because you believe it is a confounder, or exclude it because you believe it is an intervening variable, is a decision you must make based on prior research and biological plausibility.

1.3 What are suppressers and how does multivariable analysis help me to deal with them?

Suppressor variables are a type of confounder. As with confounders, a suppresser is associated with the risk factor and the outcome (Figure 1.4). The difference is that on bivariate analysis there is no effect seen between the risk factor and the outcome. But when you adjust for the suppresser, the relationship between the risk factor and the outcome become significant.

Identifying and adjusting for suppressers can lead to important findings. For example, it was unknown whether taking antiretroviral treatment would prevent HIV seroconversion among health care workers who sustained a needle stick from a patient who was HIV-infected. For several years, health care workers who had an exposure were offered zidovudine treatment, but they were told that there was no efficacy data to support its use. A randomized controlled trial was attempted, but it was disbanded because health care workers did not wish to be randomized.

✓ TIP

Statistics cannot distinguish between a confounder and an intervening variable.

✓ TIP

Unlike a typical confounder, when you have a suppresser you won't see any bivariate association between the risk factor and the outcome until you adjust for the suppresser.

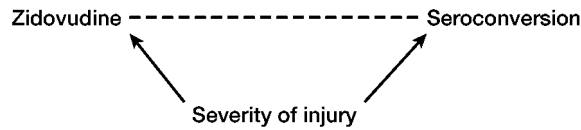


Figure 1.5. Bivariate relationships among zidovudine, severity of injury, and seroconversion.

Since a randomized controlled trial was not possible, a case-control study was performed instead.⁸ The cases were health care workers who sustained a needle stick and had seroconverted. The controls were health care workers who sustained a needle stick but had remained HIV-negative. The question was whether the proportion of persons taking zidovudine would be lower in the group who had seroconverted (the cases) than in the group who had not become infected (the controls). The investigators found that the proportion of cases using zidovudine was lower (9 of 33 cases or 27%) than the proportion of controls using zidovudine (247 of 679 controls or 36%), but the difference was not statistically significant (probability (P) = .35). Consistent with this nonsignificant trend, the odds ratio shows that zidovudine was protective (0.7), but the 95% confidence intervals were wide and did not exclude one (0.3–1.4).

However, it was known that health care workers who sustained an especially serious exposure (e.g., a deep injury or who stuck themselves with a needle that had visible blood on it) were more likely to choose to take zidovudine than health care workers who had more minor exposures. Also, health care workers who had serious exposures were more likely to seroconvert.

When the researchers adjusted their analysis for severity of injury using multiple logistic regression, zidovudine use was associated with a significantly lower risk of seroconversion (odds ratio (OR) = 0.2; 95% confidence interval (CI) = 0.1–0.6; $p < 0.01$). Thus, we have an example of a suppresser effect as shown in Figure 1.5. Severity of exposure is associated with zidovudine use and causally related to seroconversion. Zidovudine use is not associated with seroconversion in bivariate analyses but becomes significant when you adjust for severity of injury.

⁸ Cardo, D.M., Culver, D.H., Ciesielski, C.A., et al. "A case-control study of HIV seroconversion in health-care workers after percutaneous exposure." *N. Engl. J. Med.* 1997;337:1485–90.

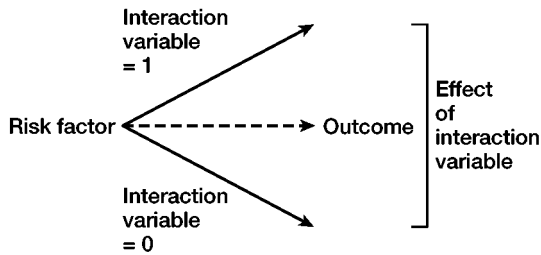


Figure 1.6. Illustration of an interaction effect.

Although this multivariable analysis demonstrated the efficacy of zidovudine on seroconversion by incorporating the suppresser variables, it should be remembered that multivariable analysis cannot adjust for other potential biases in the analysis. For example, the cases and controls for this study were not chosen from the same population, raising the possibility that selection bias may have influenced the results. Nonetheless, on the strength of this study, postexposure prophylaxis with antiretroviral treatment became the standard of care for health care workers who sustained needle sticks from HIV-contaminated needles.

1.4 What are interactions and how does multivariable analysis help me to deal with them?

An interaction occurs when the impact of a risk factor on outcome is changed by the value of a third variable. Interaction is sometimes referred to as effect modification, since the effect of the risk factor on outcome is modified by another variable.

An interaction is illustrated in Figure 1.6. The risk factor's effect on outcome (solid lines) differs depending on the value of the interaction variable (whether it is 1 or 0). The dotted line indicates the relationship without consideration of the interaction effect.

In extreme cases, an interaction may completely reverse the relationship between the risk factor and the outcome. This would occur when the risk factor increased the likelihood of outcome at one value of the interaction variable but decreased the likelihood of outcome at a different value of the interaction variable. More commonly, the effect of the risk factor on the outcome is stronger (or weaker) at certain values of the third variable.

As with confounding, stratification can be used to identify an interaction. By stratifying by the interaction variable, you can observe

DEFINITION

An *interaction* occurs when the impact of a risk factor on outcome is changed by the value of a third variable.

the effect of a risk factor on outcome at the different values of the interaction variable. You can statistically test whether the association between a risk factor and an outcome at different levels of the interaction variable are statistically different from one another using a chi-square test for homogeneity.

However, as with the use of stratification to eliminate confounding, use of stratification to demonstrate interaction has limitations. It is cumbersome to stratify by more than one or two variables; yet you may have multiple interactions in your data. Whereas stratification will accurately quantify the effect of the risk factor on the outcome at different levels of the interaction variable, this analysis will not be adjusted for the other variables in your model (e.g., confounders) that may affect the relationship between risk factor and outcome. Multi-variable analysis allows you to include interaction terms and assess them while adjusting for other variables.

For example, Zucker and colleagues evaluated whether specific signs or symptoms of myocardial infarction were different in men than in women presenting to the emergency department with chest pain or other symptoms of acute cardiac ischemia.⁹

In Table 1.7 you can see the association between the independent variables and confirmed diagnoses of acute myocardial infarction. The coefficients and odds ratios are from a multiple logistic regression model. The authors found three significant interactions involving gender: male gender and ST elevation (on electrocardiogram), male gender and congestive heart failure, and male gender and white race.

What do these interactions mean? Let's use the interaction involving male gender and ST elevations as an example (I have put these two variables and their interaction term in bold print). Note that men were more likely than women to have cardiac ischemia ($OR = 1.6$), even after adjusting for other variables associated with ischemia. Similarly, ST elevations were more likely to indicate ischemia ($OR = 8.1$). Given this, you would expect that males with ST elevations would have markedly higher risk of myocardial infarction ($1.6 \times 8.1 = 13.0$) than women ($1.0 \times 8.1 = 8.1$) (the wonderful property of odds ratios that allows you to multiply them this way is explained in Section 10.8).

The multiplication of the odds ratios of gender and ST elevations would lead you to believe that men with ST elevations would have

⁹ Zucker, D.R., Griffith, J.L., Beshansky, J.R., Selker, H.P. "Presentations of acute myocardial infarction in men and women." *J. Gen. Intern. Med.* 1997;12:79–87.

TABLE 1.7

Association of independent variables with confirmed diagnosis of acute myocardial infarction based on multiple logistic regression model.

Independent variables	Coefficients	Odds ratio
Male gender	0.4852	1.6
Age <50	0.1432	1.2
Chest pain	0.8792	2.4
Chief complaint: chest pain	0.4399	1.6
Nausea/vomiting	0.5153	1.7
Congestive heart failure	0.6759	2.0
White race	0.0987	1.1
ST elevation	2.0948	8.1
ST depression	1.2632	3.5
Q waves	0.5311	1.7
History of diabetes mellitus	0.2781	1.3
History of hypertension	0.2032	1.2
History of angina	−0.2976	0.7
History of peptic ulcers	−0.3210	0.7
Dizziness	−0.4437	0.6
Interactions		
Male gender and congestive heart failure	−0.6899	0.5
Male gender and ST elevation	−0.5187	0.6
Male gender and white race	0.5206	1.7

Adapted with permission from Zucker, D.R., et al. "Presentation of acute myocardial infarction in men and women." *J. Gen. Intern. Med.* 1997;12:79–87.

significantly higher risk of heart attack than women (13.0 vs. 8.1). In fact, the risk for men and women with ST elevations was similar. This is reflected in the negative coefficient for male gender \times ST elevations and the odds ratio of 0.6. If you multiply out the odds ratio for the interaction of male gender with ST elevations, men with ST elevations ($1.6 \times 8.1 \times 0.6 = 7.8$) and women with ST elevations ($1.0 \times 8.1 \times 1.0 = 8.1$) have a similar risk of myocardial infarction.

ST elevations are highly specific for (although not diagnostic of) myocardial infarction. It is not surprising, therefore, that the risks of

myocardial infarction are similar in men and women with ST elevations. Had being male made it even worse to have ST elevations the coefficient would have been positive, the odds ratio would have been greater than one, and we would have seen an even greater difference between the risk of heart attack for men and for women in the presence of ST elevations than the difference between 13.0 and 8.1.

Because interaction effects can be difficult to assess and interpret, I will return to this topic in Sections 8.3, 9.9, and 10.8.